

▼ Text Classification:

Data

1. we have total of 20 types of documents(Text files) and total 18828 documents(text files)
2. You can download data from this [link](#), in that you will get documents.rar folder.
If you unzip that, you will get total of 18828 documents. document name is defined as 'Class' so from document name, you can extract the label for that document.
4. Now our problem is to classify all the documents into any one of the class.
5. Below we provided count plot of all the labels in our data.

```
from google.colab import drive          #mounting the google drive
drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mour

```
!pip install bs4
!pip install tensorflow_addons
```

```
Requirement already satisfied: bs4 in /usr/local/lib/python3.7/dist-packages (0.0.1)
Requirement already satisfied: beautifulsoup4 in /usr/local/lib/python3.7/dist-packages
Collecting tensorflow_addons
  Downloading tensorflow_addons-0.14.0-cp37-cp37m-manylinux_2_12_x86_64.manylinux2010_x86_64.whl (1.1 MB)
     |██████████| 1.1 MB 12.3 MB/s
Requirement already satisfied: typeguard>=2.7 in /usr/local/lib/python3.7/dist-packages
Installing collected packages: tensorflow-addons
Successfully installed tensorflow-addons-0.14.0
```

```
import re
import os
import pandas as pd
import nltk
from nltk import ne_chunk, pos_tag, word_tokenize, sent_tokenize
from nltk.tree import Tree
nltk.download('punkt')
nltk.download('averaged_perceptron_tagger')
nltk.download('maxent_ne_chunker')
nltk.download('words')
from nltk.chunk import tree2conlltags
from sklearn.model_selection import train_test_split
from bs4 import BeautifulSoup
import numpy as np
```

```
import tensorflow as tf
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences
from tensorflow.keras.models import Model
from tensorflow.keras.layers import Input,Dense,Conv1D,Flatten,Embedding,MaxPool1D,concatenate
from tensorflow.keras.callbacks import ModelCheckpoint,TensorBoard,EarlyStopping
from tensorflow.keras.optimizers import Adam
from sklearn.preprocessing import LabelEncoder
from keras.utils import np_utils
import pickle
from tensorflow.keras.layers import Conv1D,Flatten,Embedding,MaxPool1D,concatenate,Dropout
from tensorflow.keras.callbacks import ReduceLROnPlateau

import tensorflow_addons as tfa
from tensorflow_addons.metrics import F1Score
from sklearn.metrics import f1_score
from tensorflow.keras.callbacks import EarlyStopping
from tensorflow.keras.callbacks import ModelCheckpoint
from tensorflow.keras.callbacks import ReduceLROnPlateau

[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]   /root/nltk_data...
[nltk_data]   Package averaged_perceptron_tagger is already up-to-
[nltk_data]   date!
[nltk_data] Downloading package maxent_ne_chunker to
[nltk_data]   /root/nltk_data...
[nltk_data]   Package maxent_ne_chunker is already up-to-date!
[nltk_data] Downloading package words to /root/nltk_data...
[nltk_data]   Package words is already up-to-date!

%load_ext tensorboard

### count plot of all the class labels.
```



▼ Assignment:

sci.crypt

sample document

Subject: A word of advice

From: jcopelan@nyx.cs.du.edu (The One and Only)

In article <65882@mimsy.umd.edu> mangoe@cs.umd.edu (Charley Wingate) writes:

>
>I've said 100 times that there is no "alternative" that should think you
>might have caught on by now. And there is no "alternative", but the point
>is, "rationality" isn't an alternative either. The problems of metaphysical
>and religious knowledge are unsolvable-- or I should say, humans cannot
>solve them.

How does that saying go: Those who say it can't be done shouldn't interrupt those who are doing it.

Jim

--

Have you washed your brain today?

▼ Preprocessing:

useful links: <http://www.pyregex.com/>

1. Find all emails in the document and then get the text after the "@". and then split those after that remove the words whose length is less than or equal to 2 and also remove 'com' wo

In one doc, if we have 2 or more mails, get all.

Eg:[test@dm1.d.com, test2@dm2.dm3.com]-->[dm1.d.com, dm3.dm4.com]-->[dm1,d,com,dm2,dm3,com]
append all those into one list/array. (This will give length of 18828 sentences i.e one li
Some sample output was shown below.

> In the above sample document there are emails [jcopelan@nyx.cs.du.edu, 65882@mimsy.umd.ed

preprocessing:

[jcopelan@nyx.cs.du.edu, 65882@mimsy.umd.edu, mangoe@cs.umd.edu] ==> [nyx cs du edu mimsy u
[nyx edu mimsy umd edu umd edu]

2. Replace all the emails by space in the original text.

```
< ━━━━━━ >
```

```
# we have collected all emails and preprocessed them, this is sample output
preprocessed_email
```

```
array(['juliet caltech edu',
       'coding bchs edu newsgate sps mot austlcm sps mot austlcm sps mot com dna bchs e
       'batman bmd trw', ... , 'rbdc wsnc org dscomsa desy zeus desy',
       'rbdc wsnc org morrow stanford edu pangea Stanford EDU',
       'rbdc wsnc org apollo apollo'], dtype=object)
```

```
< ━━━━━━ >
```

```
len(preprocessed_email)
```

```
18828
```

3. Get subject of the text i.e. get the total lines where "Subject:" occur and remove the word which are before the ":" remove the newlines, tabs, punctuations, any special char

Eg: if we have sentance like "Subject: Re: Gospel Dating @ \r\r\n" --> You have to get "Gos
Save all this data into another list/array.

4. After you store it in the list, Replace those sentances in original text by space.

5. Delete all the sentances where sentence starts with "Write to:" or "From:".

> In the above sample document check the 2nd line, we should remove that

6. Delete all the tags like "< anyword >"

> In the above sample document check the 4nd line, we should remove that "< 65882@mimsy.umd

7. Delete all the data which are present in the brackets.

In many text data, we observed that, they maintained the explanation of sentence

or translation of sentence to another language in brackets so remove all those.

Eg: "AAIC-The course that gets you HIRED(AAIC - Der Kurs, der Sie anstellt)" --> "AAIC-The

> In the above sample document check the 4nd line, we should remove that "(Charley Wingate)

8. Remove all the newlines('\n'), tabs('\t'), "-", "\".

9. Remove all the words which ends with ":".

Eg: "Anyword:"

> In the above sample document check the 4nd line, we should remove that "writes:"

10. Decontractions, replace words like below to full words.

please check the donors choose preprocessing for this

Eg: can't -> can not, 's -> is, i've -> i have, i'm -> i am, you're -> you are, i'll --> i

There is no order to do point 6 to 10. but you have to get final output correctly

11. Do chunking on the text you have after above preprocessing.

Text chunking, also referred to as shallow parsing, is a task that follows Part-Of-Speech Tagging and that adds more structure to the sentence.

So it combines the some phrases, named entities into single word.

So after that combine all those phrases/named entities by separating "_".

And remove the phrases/named entities if that is a "Person".

You can use `nltk.ne_chunk` to get these.

Below we have given one example. please go through it.

useful links:

<https://www.nltk.org/book/ch07.html>

<https://stackoverflow.com/a/31837224/4084039>

<http://www.nltk.org/howto/tree.html>

<https://stackoverflow.com/a/44294377/4084039>



```
#i am living in the New York
print("i am living in the New York -->", list(chunks))
print(" ")
print("-"*50)
print(" ")
#My name is Srikanth Varma
print("My name is Srikanth Varma -->", list(chunks1))
```

i am living in the New York --> [('i', 'NN'), ('am', 'VBP'), ('living', 'VBG'), ('in',

```
My name is Srikanth Varma --> [('My', 'PRP$'), ('name', 'NN'), ('is', 'VBZ'), Tree('PERS')
```

We did chunking for above two lines and then We got one list where each word is mapped to a POS(parts of speech) and also if you see "New York" and "Srikanth Varma", they got combined and represented as a tree and "New York" was referred as "GPE" and "Srika so now you have to Combine the "New York" with "_" i.e "New_York" and remove the "Srikanth Varma" from the above sentence because it is a person.

13. Replace all the digits with space i.e delete all the digits.

> In the above sample document, the 6th line have digit 100, so we have to remove that.

14. After doing above points, we observed there might be few word's like

"word" (i.e starting and ending with the _), "word" (i.e starting with the _), "word" (i.e ending with the _) remove the _ from these type of words.

15. We also observed some words like "OneLetter_word"- eg: d_berlin,

"TwoLetters_word" - eg: dr_berlin , in these words we remove the "OneLetter_" (d_berlin ==> "TwoLetters_" (de_berlin ==> berlin). i.e remove the words which are length less than or equal to 2 after splitting those words by "_".

16. Convert all the words into lower case and lowe case

and remove the words which are greater than or equal to 15 or less than or equal to 2.

17. replace all the words except "A-Za-z_" with space.

18. Now You got Preprocessed Text, email, subject. create a dataframe with those.

Below are the columns of the df.

```
data.columns
```

```
Index(['text', 'class', 'preprocessed_text', 'preprocessed_subject',
       'preprocessed_emails'],
      dtype='object')
```

```
data.iloc[400]
```

text	From: arc1@ukc.ac.uk (Tony Curtis)\r\r\r\nSubj...
------	---

```
class                                     alt.atheism
preprocessed_text      said re is article if followed the quoting rig...
preprocessed_subject                         christian morality is
preprocessed_emails                           ukc mac macalstr edu
Name: 567, dtype: object
```

```
#Extracting all documents in drive
os.chdir("/content/drive/MyDrive/Assignment19_Document_Classification_CNN/")
print(os.listdir())
```

```
[]
```

```
!unrar e documents.rar /Documents/
```

Streaming output truncated to the last 5000 lines.

Extracting /Documents/sci.space_60893.txt	OK
Extracting /Documents/sci.space_60894.txt	OK
Extracting /Documents/sci.space_60895.txt	OK
Extracting /Documents/sci.space_60896.txt	OK
Extracting /Documents/sci.space_60897.txt	OK
Extracting /Documents/sci.space_60898.txt	OK
Extracting /Documents/sci.space_60899.txt	OK
Extracting /Documents/sci.space_60900.txt	OK
Extracting /Documents/sci.space_60901.txt	OK
Extracting /Documents/sci.space_60902.txt	OK
Extracting /Documents/sci.space_60903.txt	OK
Extracting /Documents/sci.space_60904.txt	OK
Extracting /Documents/sci.space_60905.txt	OK
Extracting /Documents/sci.space_60906.txt	OK
Extracting /Documents/sci.space_60907.txt	OK
Extracting /Documents/sci.space_60908.txt	OK
Extracting /Documents/sci.space_60909.txt	OK
Extracting /Documents/sci.space_60910.txt	OK
Extracting /Documents/sci.space_60911.txt	OK
Extracting /Documents/sci.space_60912.txt	OK
Extracting /Documents/sci.space_60913.txt	OK
Extracting /Documents/sci.space_60914.txt	OK
Extracting /Documents/sci.space_60915.txt	OK
Extracting /Documents/sci.space_60916.txt	OK
Extracting /Documents/sci.space_60917.txt	OK
Extracting /Documents/sci.space_60918.txt	OK
Extracting /Documents/sci.space_60919.txt	OK
Extracting /Documents/sci.space_60920.txt	OK
Extracting /Documents/sci.space_60921.txt	OK
Extracting /Documents/sci.space_60922.txt	OK
Extracting /Documents/sci.space_60923.txt	OK
Extracting /Documents/sci.space_60924.txt	OK
Extracting /Documents/sci.space_60925.txt	OK
Extracting /Documents/sci.space_60926.txt	OK
Extracting /Documents/sci.space_60927.txt	OK
Extracting /Documents/sci.space_60928.txt	OK
Extracting /Documents/sci.space_60929.txt	OK
Extracting /Documents/sci.space_60930.txt	OK
Extracting /Documents/sci.space_60931.txt	OK
Extracting /Documents/sci.space_60932.txt	OK

Extracting	/Documents/sci.space_60933.txt	OK
Extracting	/Documents/sci.space_60934.txt	OK
Extracting	/Documents/sci.space_60935.txt	OK
Extracting	/Documents/sci.space_60936.txt	OK
Extracting	/Documents/sci.space_60937.txt	OK
Extracting	/Documents/sci.space_60938.txt	OK
Extracting	/Documents/sci.space_60939.txt	OK
Extracting	/Documents/sci.space_60940.txt	OK
Extracting	/Documents/sci.space_60941.txt	OK
Extracting	/Documents/sci.space_60942.txt	OK
Extracting	/Documents/sci.space_60943.txt	OK
Extracting	/Documents/sci.space_60944.txt	OK
Extracting	/Documents/sci.space_60945.txt	OK
Extracting	/Documents/sci.space_60946.txt	OK
Extracting	/Documents/sci.space_60947.txt	OK
Extracting	/Documents/sci.space_60948.txt	OK
Extracting	/Documents/sci.space_60949.txt	OK
Extracting	/Documents/sci.space_60950.txt	OK

To get above mentioned data frame --> Try to Write Total Preprocessing steps in One Function Named Preprocess as below.

```
def preprocess(file):                                     #preprocess function to do pr
    """Do all the Preprocessing as shown above and
    return a tuple contain preprocess_email,preprocess_subject,preprocess_text for that Text_
    dir="/Documents/"
    doc_class_label = file.split('_')[0]

    with open(dir+file,'rb') as f1:
        preprocessed_emails=[]
        preprocessed_subjects=[]
        text = f1.read()
        org_text = str(text)

        soup = BeautifulSoup(text,'lxml')                  #Getting cl
        text = soup.get_text()
        match=re.findall(r"[a-z0-9\.\\-+_]++@[a-z0-9\\.\\-+_]++\\.[a-z]+",org_text)      #rege
                                                               #Took the below reference to see h

        for mail in match:                                #loop to do basic email preprocessing like
            mail=mail.split('@')[1] #getting text after the "@". and then split those tex
            mail=mail.split('.')
            for val in mail:
                if len(val)>2 and val not in 'com': #after that remove the words whose
                    preprocessed_emails.append(val)          #appending the
            final_mail_st=""
            for mail in preprocessed_emails:
                final_mail_st+=" "+mail
            if (re.findall(r'Subject:.*',text)):
                subject = re.findall(r'Subject:.*',text)[0].split(":")[-1] ##searching for a
```

```

subject = re.sub('^[^A-Za-z0-9]+', ' ', subject)                      #searching for al
subject = subject.lower()

else:
    subject= " "                                                       #get the total l
                                                               #the word which are

text=re.sub(r"[a-z0-9\.\\-+_]@[a-z0-9\\.\\-+_]\\.[a-z]+",' ',text )      #Getting Document text

text= re.sub(r'Write to:+', ' ',text)                                     #Delete all the sentances whe
text= re.sub(r'From:+', ' ',text)

text=re.sub(r'\\([^\n]*\\)', ' ', text)                                     # Delete all the data(words s
text=re.sub(r'\\(','"',text)
text=re.sub(r'\\)',"'",text)
text=re.sub(r'\\t','"',text)                                              #Removing all the newlines('\\n'), tabs
text=re.sub(r'\\-','"',text)
text=re.sub(r'\\\\','"',text)
text=re.sub(r'\\n','"',text)
text = re.sub(r"\s+"," ",text)
text = re.sub(r"/",".",text)
text = re.sub(r"\.",","",text)
text = re.sub(r"\+","",text)
text = re.sub(r"\*","",text)
text=re.sub(r'\w+:\\s?','"',text)                                         #Removing all the words which end
#Implemeting Decontraction of words in preprocessed text
text= re.sub(r"won't", "will not",text)
text = re.sub(r"can\\'t", "can not",text)
text= re.sub(r"n\\'t", " not",text)
text = re.sub(r"\\'re", " are",text)
text= re.sub(r"\\'s", " is",text)
text = re.sub(r"\\'d", " would",text)
text= re.sub(r"\\'ll", " will",text)
text= re.sub(r"\\'t", " not",text)
text= re.sub(r"\\'ve", " have",text)
text = re.sub(r"\\'m", " am",text)

#Implementing Parts of Peech Tagging to invidual words and then doing Chunking using
chunks=ne_chunk(pos_tag(word_tokenize(text)))
Person_labels=[]
Non_Person_labels=[]
Chunked_text1=ne_chunk(pos_tag(word_tokenize(text)))
Person_labels=[]
Non_Person_labels=[]
for c in Chunked_text1:
    if isinstance(c, Tree) and c.label() == 'PERSON':                  #Checking if the ch
        # The object is <class 'nltk.tree.Tree'> and label is PERSON Entity
        Person_labels.append(' '.join([w for w, _ in c.leaves()]))
    if isinstance(c, Tree) and c.label()!='PERSON':                      #Checking if
        if len(c)>2:
            Non_Person_labels.append(' '.join([w for w, _ in c.leaves()]))

```

```

        else:
            Non_Person_labels.append(c.leaves()[0][0])
    for word in Person_labels:                                #Replacing all the PERSON(Nam
        text=re.sub(word,"",text)
    for word in Non_Person_labels:                            #Joining all the NON-PERSON
        old_word_list=word.split("_")
        org_non_person_word=" ".join(old_word_list)
        text=re.sub(org_non_person_word,word,text)

    text=re.sub(r'[0-9]+'," ",text)                         #Replacing all the digits with space
    text=re.sub(r'\.','','',text)
    text = re.sub(r"(_?)([A-Za-z0-9])(_?)",r'\2',text) #remove the _ from these type of w
    text = re.sub(r"([A-Za-z]{1,2})(_)(A-Za-z)",r"\g<3>",text) #remove the words which ar
    text=re.sub(r"^[A-Za-z_]", " ",text)                  #Finally replacing all the words except
    text=text.lower()                                     #Converting all the words to lower case in preprocessed

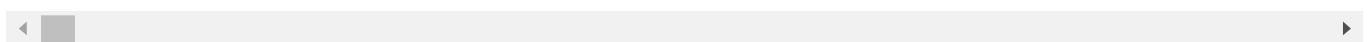
    final_preprocessed_text = ' '.join([wrds for wrds in text.split() if len(wrds)>2 and len

    return [org_text,doc_class_label,final_preprocessed_text,subject,final_mail_st]

```

```
#Running the preprocess function for sample alt.atheism_49960.txt file
dir="/Documents/" #Directory path where all documents are present
sample_dat_lst=preprocess('alt.atheism_49960.txt')
sample_dat_lst
```

```
[ 'b\'From: mathew <mathew@mantis.co.uk>\nSubject: Alt.Atheism FAQ: Atheist Resources\\r
 'alt.atheism',
 'resources december atheist organizations usa freedom from religion foundation fish bum
 ' atheist resources',
 ' mantis netcom mantis']
```



```
#Preprocessing every document.
dir="/Documents/"      #directory where all file are present in drive
f_processed=0           #counter variable
data_list=[]            #finla list to store data of each document as list
for file in os.listdir(dir):          #loop to iterate over each file
    if f_processed%1000==0:
        print("Number of files processed:",f_processed)
    if f_processed==18828:
        print("All files processed",f_processed)
```

```
f_processed+=1
final_dat=preprocess(file)
data_list.append(final_dat)                                #calling the text preprocess() function
                                                               #appending the data to final list for each fi
```

```
Number of files processed: 0
Number of files processed: 1000
Number of files processed: 2000
Number of files processed: 3000
Number of files processed: 4000
Number of files processed: 5000
Number of files processed: 6000
Number of files processed: 7000
Number of files processed: 8000
Number of files processed: 9000
Number of files processed: 10000
Number of files processed: 11000
Number of files processed: 12000
Number of files processed: 13000
Number of files processed: 14000
Number of files processed: 15000
Number of files processed: 16000
Number of files processed: 17000
Number of files processed: 18000
```

```
#Making the final dataframe
final_df1 = pd.DataFrame(data_list,columns =['text', 'class', 'preprocessed_text', 'preprocessed_mail'])
```

```
final_df1.shape      #Checking the shape of final data
```

```
(18828, 5)
```

```
final_df1.head()          #Reading the df data
```

		text	class	preprocessed_text	preprocessed_mail
0	b"From: ricardo@rchland.vnet.ibm.com (Ricardo ...		comp.graphics	article quvdoinn concerning the proposed newsg...	newscl
1	b"From: holmes7000@iscsvax.uni.edu\nSubject: A...		comp.graphics	what the best way archive gif zip them and the...	arc
2	b'From: kxgst1+@pitt.edu (Kenneth Gilbert)\nSu...		sci.med	article chronic persistent hepatitis usually d...	persistent \

```
#Saving the final dataframe as pickle file.
```

```
os.chdir("/content/drive/MyDrive/Assignment19_Document_Classification_CNN/")
final_df1.to_pickle("final_preprocessed_text_df")

os.chdir("/content/drive/MyDrive/Assignment19_Document_Classification_CNN/")
final_df2=pd.read_pickle("final_preprocessed_text_df")           #Reading the final dataframe
final_df2.shape

(18828, 5)
```

Code checking:

After Writing preprocess function. call that functoin with the input text of 'alt.atheism_49960' doc and print the output of the preprocess function
This will help us to evaluate faster, based on the output we can suggest you if there are any changes.

After writing Preprocess function, call the function for each of the document(18828 docs) and then create a dataframe as mentioned above.

Training The models to Classify:

1. Combine "preprocessed_text", "preprocessed_subject", "preprocessed_emails" into one column
2. Now Split the data into Train and test. use 25% for test also do a stratify split.
3. Analyze your text data and pad the sequence if required.
Sequence length is not restricted, you can use anything of your choice.
you need to give the reasoning
4. Do Tokenizer i.e convert text into numbers. please be careful while doing it.
if you are using tf.keras "Tokenizer" API, it removes the "_", but we need that.
5. code the model's (Model-1, Model-2) as discussed below
and try to optimize that models.
6. For every model use predefined Glove vectors.
Don't train any word vectors while Training the model.
7. Use "categorical_crossentropy" as Loss.
8. Use **Accuracy and Micro Avggeraged F1 score** as your as Key metrics to evaluate your model.

9. Use Tensorboard to plot the loss and Metrics based on the epoches.

10. Please save your best model weights in to '**best_model_L.h5**' (L = 1 or 2).

11. You are free to choose any Activation function, learning rate, optimizer.
But have to use the same architecture which we are giving below.

12. You can add some layer to our architecture but you **deletion** of layer is not acceptable.

13. Try to use **Early Stopping** technique or any of the callback techniques that you did in t

14. For Every model save your model to image (Plot the model) with shapes
and include those images in the notebook markdown cell,
upload those images to Classroom. You can use "plot_model"
please refer [this](#) if you don't know how to plot the model with shapes.

Model-1: Using 1D convolutions with word embeddings

Encoding of the Text --> For a given text data create a Matrix with Embedding layer as shown.
In the example we have considered d = 5, but in this assignment we will get d = dimension of
i.e if we have maximum of 350 words in a sentence and embedding of 300 dim word vector,

we result in 350*300 dimensional matrix for each sentence as output after embedding layer

The diagram illustrates the process of generating word embeddings. On the left, the sentence "I like this movie very much !" is written vertically. To its right is a 7x5 grid of numbers, representing the 350-dimensional output of an embedding layer for each word in the sentence. The grid contains the following data:

0.6	0.5	0.2	-0.1	0.4
0.8	0.9	0.1	0.5	0.1
0.4	0.6	0.1	-0.1	0.7
---	---	---	---	---
---	---	---	---	---
---	---	---	---	---
---	---	---	---	---

Ref: <https://i.imgur.com/kiVQuk1.png>

Reference:

<https://stackoverflow.com/a/43399308/4084039>

<https://missinglink.ai/guides/keras/keras-conv1d-working-1d-convolutional-neural-networks-k>

How EMBEDDING LAYER WORKS

-
- Go through this blog, if you have any doubt on using predefined Embedding values in Embedding layer - <https://machinelearningmastery.com/use-word-embedding-layers-deep-learning-keras/>
- ▼ values in Embedding layer - <https://machinelearningmastery.com/use-word-embedding-layers-deep-learning-keras/>



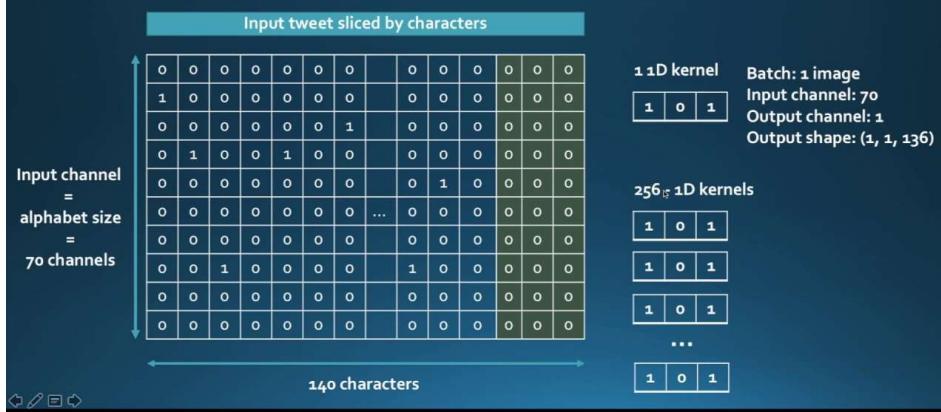
ref: '<https://i.imgur.com/fv1GvFJ.png>'

1. all are Conv1D layers with any number of filter and filter sizes, there is no restriction.
2. use concatenate layer is to concatenate all the filters/channels.
3. You can use any pool size and stride for maxpooling layer.
4. Don't use more than 16 filters in one Conv layer because it will increase the no of parameters
(Only recommendation if you have less computing power)
5. You can use any number of layers after the Flatten Layer.



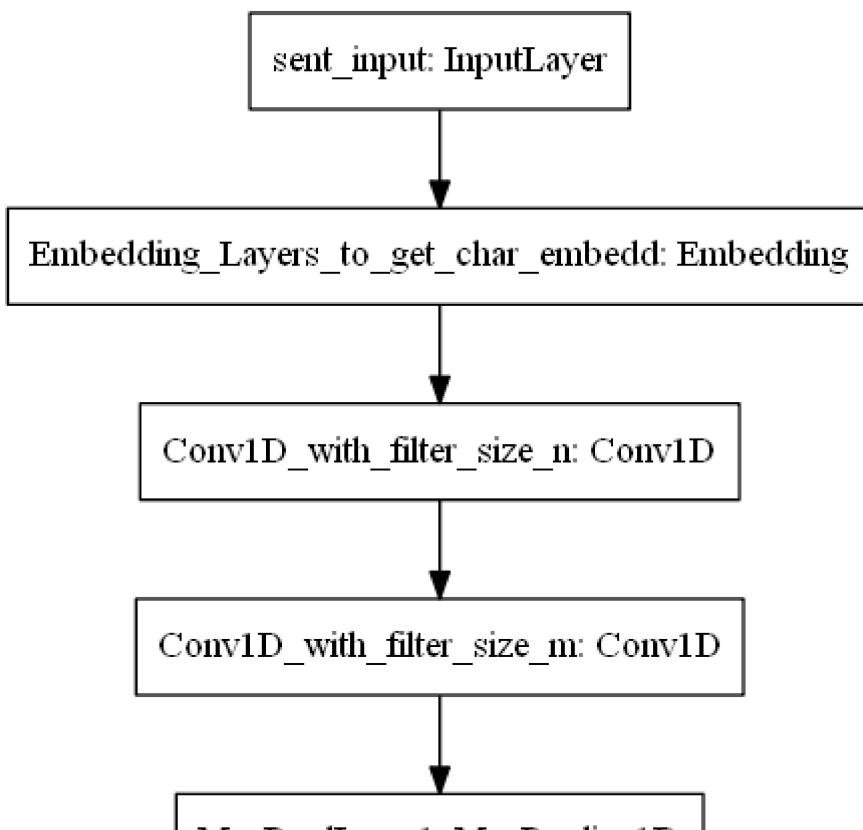
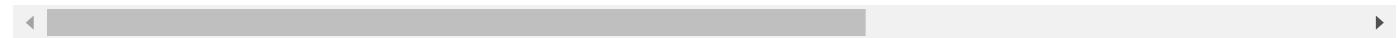
▼ Model-2 : Using 1D convolutions with character embedding

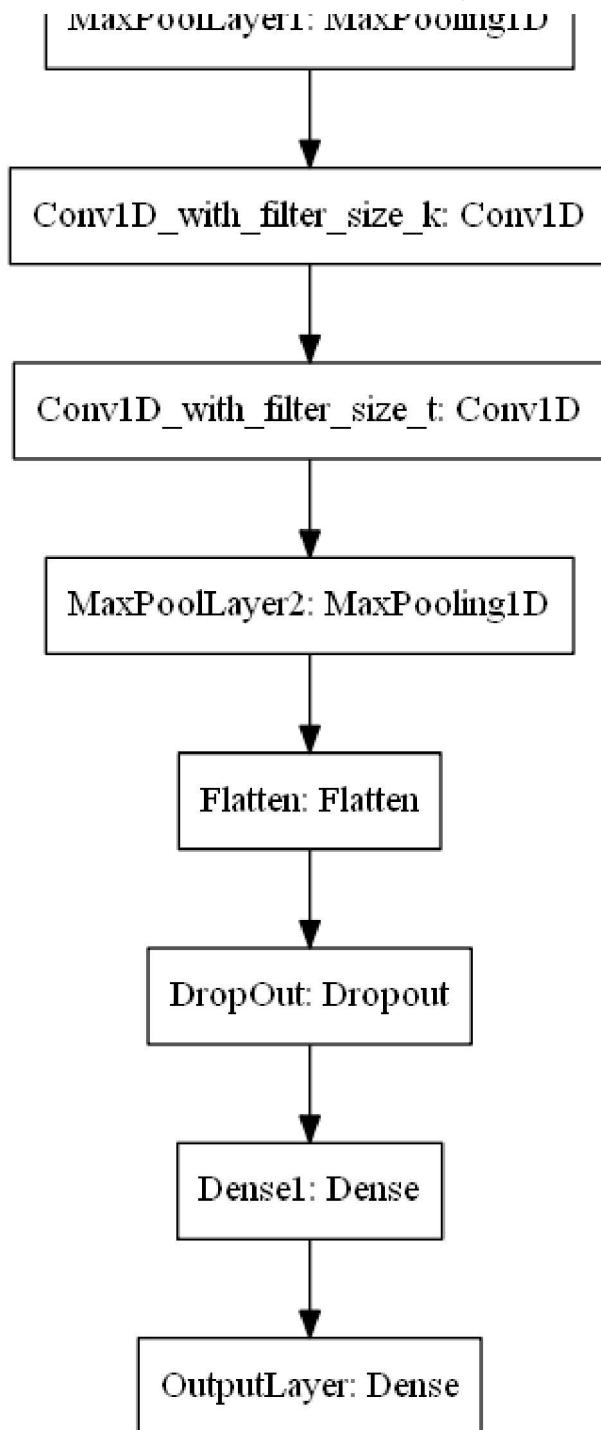
Use 1D-convolutions!



Here are the some papers based on Char-CNN

1. Xiang Zhang, Junbo Zhao, Yann LeCun. [Character-level Convolutional Networks for Text](#)
 2. Yoon Kim, Yacine Jernite, David Sontag, Alexander M. Rush. [Character-Aware Neural Lar](#)
 3. Shaojie Bai, J. Zico Kolter, Vladlen Koltun. [An Empirical Evaluation of Generic Conv](#)
 4. Use the pretrained char embeddings <https://github.com/minimaxir/char-embeddings/blob/>





```
#Combining preprocessed text ,subject and email and creating df column for same
final_df2['pre_text_pre_suject_pre_email']=final_df2['preprocessed_text']+final_df2['preprocess
final_df2.head()
```

	text	class	preprocessed_text	preprocessed
0	b"From: ricardo@rchland.vnet.ibm.com (Ricardo ...	comp.graphics	article quvdoinn concerning the proposed newsg...	news...
1	b"From: holmes7000@iscsvax.uni.edu\nSubject: A...	comp.graphics	what the best way archive gif zip them and the...	arc...

```
#checking the datatype of df columns
```

```
final_df2.dtypes
```

text	object
class	object
preprocessed_text	object
preprocessed_subject	object
preprocessed_mail	object
pre_text_pre_suject_pre_email	object
dtype:	object

```
X=final_df2['pre_text_pre_suject_pre_email']
```

```
Y=final_df2['class']
```

```
#encoding Target Y class labels
```

```
Label_en=LabelEncoder()
```

```
Label_en.fit(Y)
```

```
Label_en_Y=Label_en.transform(Y)
```

```
#Converting final labels to matrix
```

```
Y=np_utils.to_categorical(Label_en_Y)
```

```
Y.shape
```

```
(18828, 20)
```

```
#Splitting the data into train and test
```

```
X_train,X_test,Y_train,Y_test=train_test_split(X,Y,test_size=0.25,stratify=Y)
```

```
#Finding the optimal sequence length
```

```
sents_len=[len(sen) for sen in X_train]
```

```
sents_len=sorted(sents_len)
```

```
sents_len=np.array(sents_len)
```

```
#checking the stats for sequence selection
```

```
Median=int(np.percentile(sents_len,95))
```

```
Median
```

```
3112
```

```
#Checking 95th percentile
```

```
per_tile_95=int(np.percentile(sents_len,95))
```

```
per_tile_95
```

```
3112
```

```
#Checking 98th percentile
per_tile_98=int(np.percentile(sents_len,98))
per_tile_98
```

```
5346
```

```
max_length_of_sent=per_tile_98
max_length_of_sent
```

```
5346
```

Since 98 Percent of text sentences are less than 5241 .so we
 will choose max len=5241

```
ke_tokenizer = Tokenizer(filters='!"#$%&()*+,-./:;=>?@[\\]^`{|}~\t\n')    # Referred the foll
ke_tokenizer.fit_on_texts(X_train)
X_train = ke_tokenizer.texts_to_sequences(X_train)
X_test = ke_tokenizer.texts_to_sequences(X_test)
```

```
#Applying padding based on maxlen
X_train=pad_sequences(X_train,max_length_of_sent,padding='post')
X_test=pad_sequences(X_test,max_length_of_sent,padding='post')
```

```
print(X_train.shape)
print(X_test.shape)
```

```
(14121, 5346)
(4707, 5346)
```

```
Y.shape
```

```
(18828, 20)
```

```
!wget http://nlp.stanford.edu/data/glove.6B.zip
```

```
--2021-11-01 05:53:27-- http://nlp.stanford.edu/data/glove.6B.zip
Resolving nlp.stanford.edu (nlp.stanford.edu)... 171.64.67.140
Connecting to nlp.stanford.edu (nlp.stanford.edu)|171.64.67.140|:80... connected.
HTTP request sent, awaiting response... 302 Found
Location: https://nlp.stanford.edu/data/glove.6B.zip [following]
```

```
--2021-11-01 05:53:27-- https://nlp.stanford.edu/data/glove.6B.zip
Connecting to nlp.stanford.edu (nlp.stanford.edu)|171.64.67.140|:443... connected.
HTTP request sent, awaiting response... 301 Moved Permanently
Location: http://downloads.cs.stanford.edu/nlp/data/glove.6B.zip [following]
--2021-11-01 05:53:28-- http://downloads.cs.stanford.edu/nlp/data/glove.6B.zip
Resolving downloads.cs.stanford.edu (downloads.cs.stanford.edu)... 171.64.64.22
Connecting to downloads.cs.stanford.edu (downloads.cs.stanford.edu)|171.64.64.22|:80...
HTTP request sent, awaiting response... 200 OK
Length: 862182613 (822M) [application/zip]
Saving to: 'glove.6B.zip'

glove.6B.zip      100%[=====] 822.24M  4.98MB/s    in 2m 41s

2021-11-01 05:56:08 (5.12 MB/s) - 'glove.6B.zip' saved [862182613/862182613]
```



```
print(os.getcwd())

/content/drive/My Drive/Assignment19_Document_Classification_CNN

!unzip glove*.zip

Archive: glove.6B.zip
  inflating: glove.6B.50d.txt
  inflating: glove.6B.100d.txt
  inflating: glove.6B.200d.txt
  inflating: glove.6B.300d.txt

#Indexing the word vectors                      #used the below reference to see how load g

embd_dict = {}
f = open('glove.6B.50d.txt', encoding='utf-8')
for line in f:
    values = line.split()
    word = values[0]
    coefs = np.asarray(values[1:], dtype='float32')
    embd_dict[word] = coefs
f.close()

print("No of word vetors ",len(embd_dict))

No of word vetors  400000

#Creating embedding Matrix
emd_mat_size = len(ke_tokenizer.word_index)+1
embd_matrix = np.zeros((emd_mat_size,50))

for word,i in ke_tokenizer.word_index.items():
```

```

emb_word = embd_dict.get(word)
if emb_word is not None:
    embd_matrix[i]=emb_word

print("Embedding Matrix shape:",embd_matrix.shape)

```

Embedding Matrix shape: (83617, 50)

▼ 1st Model

```

#Defining Embedding Layer
Embd_ly=Embedding(len(ke_tokenizer.word_index)+1,50,embeddings_initializer=tf.keras.initializ

input_layer=Input(shape=(max_length_of_sent,))
embd_layer=Embd_ly(input_layer)
M1=Conv1D(64,4,activation='relu',kernel_initializer=tf.keras.initializers.he_normal(),kernel_
N1=Conv1D(64,4,activation='relu',kernel_initializer=tf.keras.initializers.he_normal(),kernel_
O1=Conv1D(64,4,activation='relu',kernel_initializer=tf.keras.initializers.he_normal(),kernel_
sec_layer=concatenate([M1,N1,O1])
max_pool1=MaxPool1D(3)(sec_layer)
I1=Conv1D(32,8,activation='relu',kernel_initializer=tf.keras.initializers.he_normal(),kernel_
J1=Conv1D(32,8,activation='relu',kernel_initializer=tf.keras.initializers.he_normal(),kernel_
K1=Conv1D(32,8,activation='relu',kernel_initializer=tf.keras.initializers.he_normal(),kernel_
third_layer=concatenate([I1,J1,K1])
max_pool2=MaxPool1D(3)(third_layer)
fourth_layer=Conv1D(32,4,activation='relu',kernel_initializer=tf.keras.initializers.he_normal
flat_layer1=Flatten()(fourth_layer)
dropout_layer1=Dropout(0.2)(flat_layer1)
dense_layer1=Dense(64,activation='relu',kernel_initializer=tf.keras.initializers.he_normal(),
Output_layer1=Dense(20,activation='softmax',kernel_initializer=tf.keras.initializers.glorot_n
M1=Model(inputs=input_layer,outputs=Output_layer1)

#Custom callback for F1_score
class f1(tf.keras.callbacks.Callback):
    def __init__(self,validation_data):
        super().__init__()
        self.validation_data=validation_data                      #initializing the X_test,Y_test
    def on_train_begin(self,logs={}):
        self.f1_scores=[]                                     #list to store f1_scores

    def on_epoch_end(self,epoch,logs={}):
        y_predicted =self.model.predict(self.validation_data[0]).argmax(axis=1)    #getting
        y_true=np.zeros(Y_test.shape[0])
        for x in range(len(y_true)):
            y_true[x]=int(np.argmax(Y_test[x]))

```

```

val_f1=f1_score(y_true, y_predicted,average='micro') #computin f1_score

self.f1_scores.append(val_f1)
print(" Validation F1 score {}".format(val_f1))

```

M1.summary()

Model: "model_2"

Layer (type)	Output Shape	Param #	Connected to
<hr/>			
input_3 (InputLayer)	[(None, 5346)]	0	
embedding_1 (Embedding)	(None, 5346, 50)	4180850	input_3[0][0]
conv1d_14 (Conv1D)	(None, 5343, 64)	12864	embedding_1[2][0]
conv1d_15 (Conv1D)	(None, 5343, 64)	12864	embedding_1[2][0]
conv1d_16 (Conv1D)	(None, 5343, 64)	12864	embedding_1[2][0]
concatenate_4 (Concatenate)	(None, 5343, 192)	0	conv1d_14[0][0] conv1d_15[0][0] conv1d_16[0][0]
max_pooling1d_4 (MaxPooling1D)	(None, 1781, 192)	0	concatenate_4[0][0]
conv1d_17 (Conv1D)	(None, 1774, 32)	49184	max_pooling1d_4[0][0]
conv1d_18 (Conv1D)	(None, 1774, 32)	49184	max_pooling1d_4[0][0]
conv1d_19 (Conv1D)	(None, 1774, 32)	49184	max_pooling1d_4[0][0]
concatenate_5 (Concatenate)	(None, 1774, 96)	0	conv1d_17[0][0] conv1d_18[0][0] conv1d_19[0][0]
max_pooling1d_5 (MaxPooling1D)	(None, 591, 96)	0	concatenate_5[0][0]
conv1d_20 (Conv1D)	(None, 588, 32)	12320	max_pooling1d_5[0][0]
flatten_2 (Flatten)	(None, 18816)	0	conv1d_20[0][0]
dropout_2 (Dropout)	(None, 18816)	0	flatten_2[0][0]
dense_4 (Dense)	(None, 64)	1204288	dropout_2[0][0]
dense_5 (Dense)	(None, 20)	1300	dense_4[0][0]
<hr/>			
Total params: 5,584,902			
Trainable params: 1,404,052			
Non-trainable params: 4,180,850			



```
## f1_score_callback
f1_metric=f1((X_test,Y_test))

## Callback for saving best model
checkpoint = ModelCheckpoint(filepath='best_model_1.h5',verbose=1,monitor='val_accuracy',
                             mode='max',save_best_only=True)

## Callback for earlystopping
early_stop = EarlyStopping(monitor="val_accuracy",mode='max',patience=2)

## Tensorboard
log_dir = "logs"
tensorboard_cb = TensorBoard(log_dir=log_dir,histogram_freq=1,write_graph=True)

## all callbacks
callbacks =[f1_metric,early_stop,tensorboard_cb,checkpoint]

## compile model
M1.compile(loss='categorical_crossentropy', optimizer=Adam(learning_rate=0.001), metrics=['ac
```

Trainning

```
M1.fit(X_train,Y_train,epochs=20,verbose=2,validation_data=(X_test,Y_test),batch_size =128,ca
      Validation F1 score 0.36520076481835567

Epoch 00003: val_accuracy improved from 0.34608 to 0.36520, saving model to best_mode
Epoch 4/20
111/111 - 547s - loss: 2.2386 - accuracy: 0.4127 - val_loss: 2.2070 - val_accuracy: 0
      Validation F1 score 0.4280858296154663

Epoch 00004: val_accuracy improved from 0.36520 to 0.42809, saving model to best_mode
Epoch 5/20
111/111 - 546s - loss: 2.1249 - accuracy: 0.4441 - val_loss: 2.0940 - val_accuracy: 0
      Validation F1 score 0.46483960059485874

Epoch 00005: val_accuracy improved from 0.42809 to 0.46484, saving model to best_mode
Epoch 6/20
111/111 - 547s - loss: 2.0690 - accuracy: 0.4601 - val_loss: 2.0336 - val_accuracy: 0
      Validation F1 score 0.47142553643509666

Epoch 00006: val_accuracy improved from 0.46484 to 0.47143, saving model to best_mode
Epoch 7/20
111/111 - 544s - loss: 2.0210 - accuracy: 0.4784 - val_loss: 1.9575 - val_accuracy: 0
      Validation F1 score 0.4981941788825154

Epoch 00007: val_accuracy improved from 0.47143 to 0.49819, saving model to best_mode
Epoch 8/20
111/111 - 548s - loss: 1.9754 - accuracy: 0.4917 - val_loss: 1.9641 - val_accuracy: 0
      Validation F1 score 0.5052050138092203
```

```
Epoch 00008: val_accuracy improved from 0.49819 to 0.50521, saving model to best_mode
Epoch 9/20
111/111 - 547s - loss: 1.9373 - accuracy: 0.5063 - val_loss: 1.9129 - val_accuracy: 0
Validation F1 score 0.5272997663055025

Epoch 00009: val_accuracy improved from 0.50521 to 0.52730, saving model to best_mode
Epoch 10/20
111/111 - 546s - loss: 1.9314 - accuracy: 0.5085 - val_loss: 1.8865 - val_accuracy: 0
Validation F1 score 0.5400467388995114

Epoch 00010: val_accuracy improved from 0.52730 to 0.54005, saving model to best_mode
Epoch 11/20
111/111 - 545s - loss: 1.8793 - accuracy: 0.5271 - val_loss: 1.8821 - val_accuracy: 0
Validation F1 score 0.5481198215423837

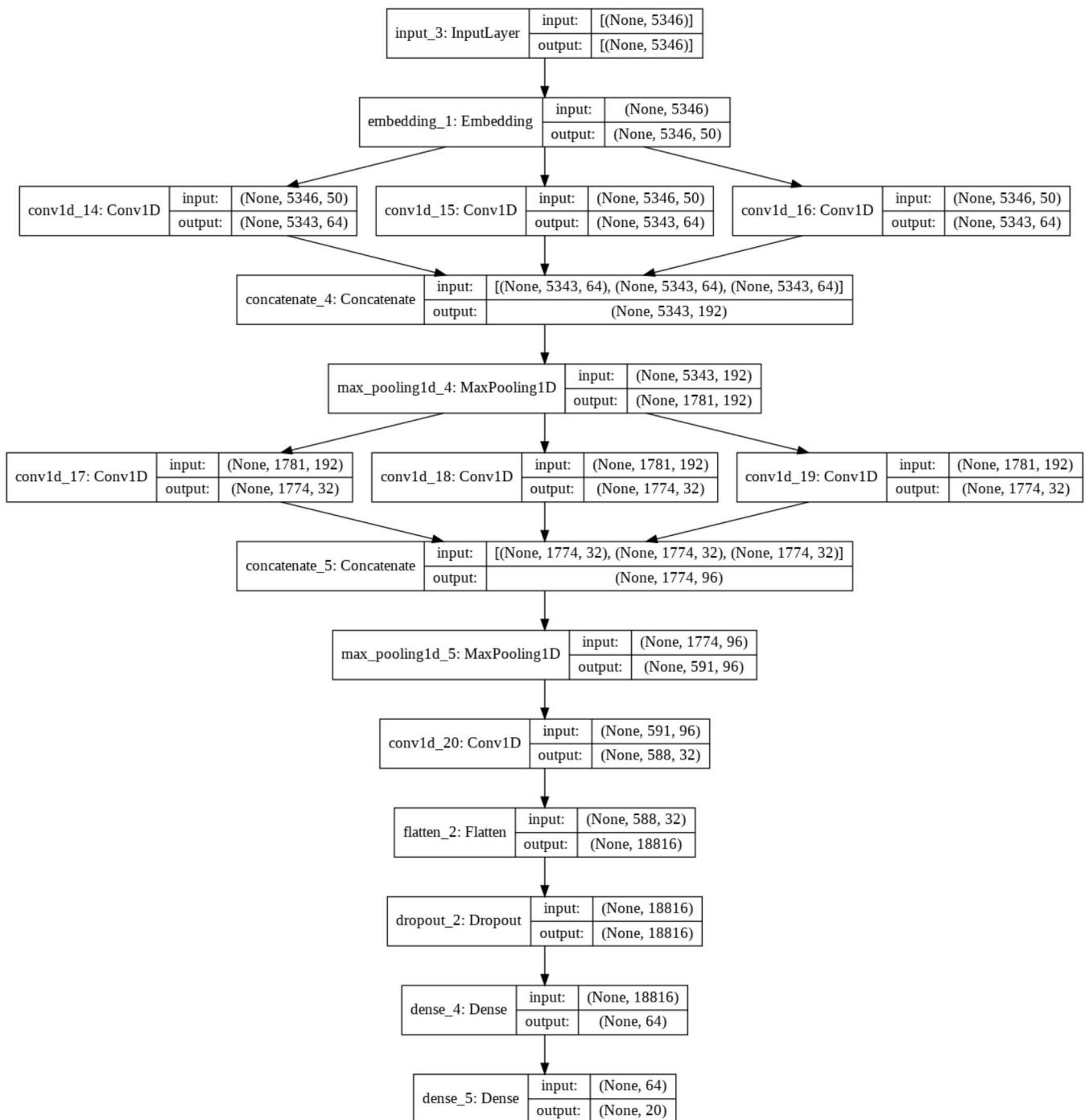
Epoch 00011: val_accuracy improved from 0.54005 to 0.54812, saving model to best_mode
Epoch 12/20
111/111 - 544s - loss: 1.8790 - accuracy: 0.5320 - val_loss: 1.8691 - val_accuracy: 0
Validation F1 score 0.5417463352453792

Epoch 00012: val_accuracy did not improve from 0.54812
Epoch 13/20
111/111 - 546s - loss: 1.8437 - accuracy: 0.5422 - val_loss: 1.8618 - val_accuracy: 0
Validation F1 score 0.5593796473337582

Epoch 00013: val_accuracy improved from 0.54812 to 0.55938, saving model to best_mode
Epoch 14/20
111/111 - 545s - loss: 1.8536 - accuracy: 0.5519 - val_loss: 1.8788 - val_accuracy: 0
Validation F1 score 0.5398342893562779

Epoch 00014: val_accuracy did not improve from 0.55938
```

```
tf.keras.utils.plot_model(M1, to_file='M1.png', show_shapes=True, show_layer_names=True)
```



```
%tensorboard --logdir logs
```

TensorBoard

SCALARS

GRAPHS

INACTIVE

- Show data download links
- Ignore outliers in chart scaling

Tooltip sorting method: default

Smoothing



0.6

Horizontal Axis

STEP RELATIVE

WALL

Runs

Write a regex to filter runs

train

validation

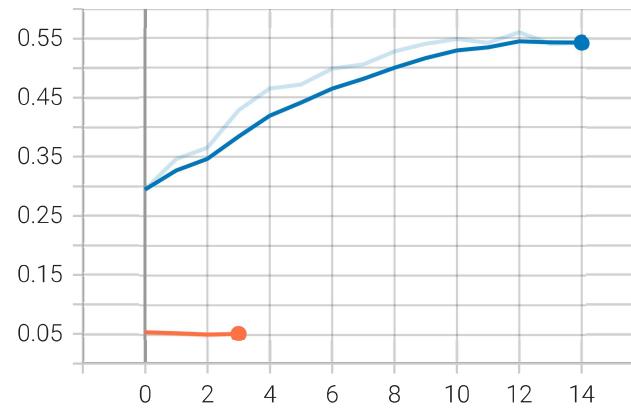
TOGGLE ALL RUNS

logs

Filter tags (regular expressions supported)

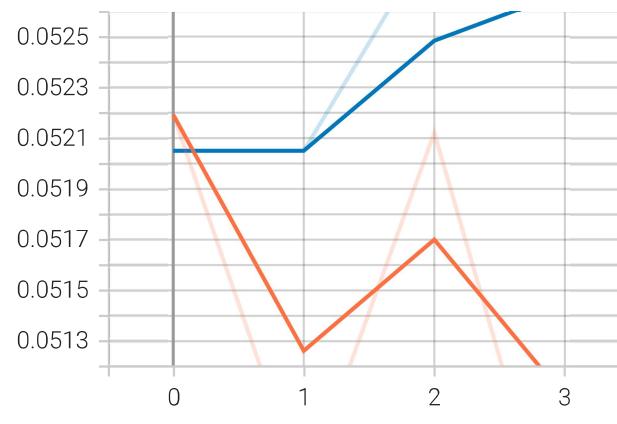
epoch_accuracy

epoch_accuracy
tag: epoch_accuracy



epoch_f1_score

epoch_f1_score
tag: epoch_f1_score



Model2

```
#Splitting the data into train and test
X_train,X_test,Y_train,Y_test=train_test_split(X,Y,test_size=0.25,stratify=Y)
```

```
ke_tokenizer2= Tokenizer(filters='!"#$%&()*+,-./:;=>?@[\\]^`{|}~\t\n',char_level=True,oov_to

ke_tokenizer2.fit_on_texts(X_train)

print(ke_tokenizer2.word_index)      #Getting the word vocab of chars

{'UNK': 1, ' ': 2, 'e': 3, 't': 4, 'a': 5, 'o': 6, 'n': 7, 'i': 8, 's': 9, 'r': 10, 'h': 11, 'd': 12, 'l': 13, 'c': 14, 'u': 15, 'g': 16, 'm': 17, 'y': 18, 'p': 19, 'b': 20, 'f': 21, 'j': 22, 'h': 23, 'k': 24, 'v': 25, 'w': 26, 'x': 27, 'z': 28, 'q': 29, 's': 30, 't': 31, 'n': 32, 'r': 33, 'l': 34, 'o': 35, 'u': 36, 'i': 37, 'e': 38, 'a': 39, 'm': 40, 'h': 41, 'd': 42, 'c': 43, 'g': 44, 'p': 45, 'f': 46, 'j': 47, 'v': 48, 'w': 49, 'x': 50, 'z': 51, 'q': 52, 's': 53, 't': 54, 'n': 55, 'r': 56, 'l': 57, 'o': 58, 'u': 59, 'i': 60, 'e': 61, 'a': 62, 'm': 63, 'h': 64, 'd': 65, 'c': 66, 'g': 67, 'p': 68, 'f': 69, 'j': 70, 'v': 71, 'w': 72, 'x': 73, 'z': 74, 'q': 75, 's': 76, 't': 77, 'n': 78, 'r': 79, 'l': 80, 'o': 81, 'u': 82, 'i': 83, 'e': 84, 'a': 85, 'm': 86, 'h': 87, 'd': 88, 'c': 89, 'g': 90, 'p': 91, 'f': 92, 'j': 93, 'v': 94, 'w': 95, 'x': 96, 'z': 97, 'q': 98, 's': 99, 't': 100, 'n': 101, 'r': 102, 'l': 103, 'o': 104, 'u': 105, 'i': 106, 'e': 107, 'a': 108, 'm': 109, 'h': 110, 'd': 111, 'c': 112, 'g': 113, 'p': 114, 'f': 115, 'j': 116, 'v': 117, 'w': 118, 'x': 119, 'z': 120, 'q': 121, 's': 122, 't': 123, 'n': 124, 'r': 125, 'l': 126, 'o': 127, 'u': 128, 'i': 129, 'e': 130, 'a': 131, 'm': 132, 'h': 133, 'd': 134, 'c': 135, 'g': 136, 'p': 137, 'f': 138, 'j': 139, 'v': 140, 'w': 141, 'x': 142, 'z': 143, 'q': 144, 's': 145, 't': 146, 'n': 147, 'r': 148, 'l': 149, 'o': 150, 'u': 151, 'i': 152, 'e': 153, 'a': 154, 'm': 155, 'h': 156, 'd': 157, 'c': 158, 'g': 159, 'p': 160, 'f': 161, 'j': 162, 'v': 163, 'w': 164, 'x': 165, 'z': 166, 'q': 167, 's': 168, 't': 169, 'n': 170, 'r': 171, 'l': 172, 'o': 173, 'u': 174, 'i': 175, 'e': 176, 'a': 177, 'm': 178, 'h': 179, 'd': 180, 'c': 181, 'g': 182, 'p': 183, 'f': 184, 'j': 185, 'v': 186, 'w': 187, 'x': 188, 'z': 189, 'q': 190, 's': 191, 't': 192, 'n': 193, 'r': 194, 'l': 195, 'o': 196, 'u': 197, 'i': 198, 'e': 199, 'a': 200, 'm': 201, 'h': 202, 'd': 203, 'c': 204, 'g': 205, 'p': 206, 'f': 207, 'j': 208, 'v': 209, 'w': 210, 'x': 211, 'z': 212, 'q': 213, 's': 214, 't': 215, 'n': 216, 'r': 217, 'l': 218, 'o': 219, 'u': 220, 'i': 221, 'e': 222, 'a': 223, 'm': 224, 'h': 225, 'd': 226, 'c': 227, 'g': 228, 'p': 229, 'f': 230, 'j': 231, 'v': 232, 'w': 233, 'x': 234, 'z': 235, 'q': 236, 's': 237, 't': 238, 'n': 239, 'r': 240, 'l': 241, 'o': 242, 'u': 243, 'i': 244, 'e': 245, 'a': 246, 'm': 247, 'h': 248, 'd': 249, 'c': 250, 'g': 251, 'p': 252, 'f': 253, 'j': 254, 'v': 255, 'w': 256, 'x': 257, 'z': 258, 'q': 259, 's': 260, 't': 261, 'n': 262, 'r': 263, 'l': 264, 'o': 265, 'u': 266, 'i': 267, 'e': 268, 'a': 269, 'm': 270, 'h': 271, 'd': 272, 'c': 273, 'g': 274, 'p': 275, 'f': 276, 'j': 277, 'v': 278, 'w': 279, 'x': 280, 'z': 281, 'q': 282, 's': 283, 't': 284, 'n': 285, 'r': 286, 'l': 287, 'o': 288, 'u': 289, 'i': 290, 'e': 291, 'a': 292, 'm': 293, 'h': 294, 'd': 295, 'c': 296, 'g': 297, 'p': 298, 'f': 299, 'j': 300, 'v': 301, 'w': 302, 'x': 303, 'z': 304, 'q': 305, 's': 306, 't': 307, 'n': 308, 'r': 309, 'l': 310, 'o': 311, 'u': 312, 'i': 313, 'e': 314, 'a': 315, 'm': 316, 'h': 317, 'd': 318, 'c': 319, 'g': 320, 'p': 321, 'f': 322, 'j': 323, 'v': 324, 'w': 325, 'x': 326, 'z': 327, 'q': 328, 's': 329, 't': 330, 'n': 331, 'r': 332, 'l': 333, 'o': 334, 'u': 335, 'i': 336, 'e': 337, 'a': 338, 'm': 339, 'h': 340, 'd': 341, 'c': 342, 'g': 343, 'p': 344, 'f': 345, 'j': 346, 'v': 347, 'w': 348, 'x': 349, 'z': 350, 'q': 351, 's': 352, 't': 353, 'n': 354, 'r': 355, 'l': 356, 'o': 357, 'u': 358, 'i': 359, 'e': 360, 'a': 361, 'm': 362, 'h': 363, 'd': 364, 'c': 365, 'g': 366, 'p': 367, 'f': 368, 'j': 369, 'v': 370, 'w': 371, 'x': 372, 'z': 373, 'q': 374, 's': 375, 't': 376, 'n': 377, 'r': 378, 'l': 379, 'o': 380, 'u': 381, 'i': 382, 'e': 383, 'a': 384, 'm': 385, 'h': 386, 'd': 387, 'c': 388, 'g': 389, 'p': 390, 'f': 391, 'j': 392, 'v': 393, 'w': 394, 'x': 395, 'z': 396, 'q': 397, 's': 398, 't': 399, 'n': 400, 'r': 401, 'l': 402, 'o': 403, 'u': 404, 'i': 405, 'e': 406, 'a': 407, 'm': 408, 'h': 409, 'd': 410, 'c': 411, 'g': 412, 'p': 413, 'f': 414, 'j': 415, 'v': 416, 'w': 417, 'x': 418, 'z': 419, 'q': 420, 's': 421, 't': 422, 'n': 423, 'r': 424, 'l': 425, 'o': 426, 'u': 427, 'i': 428, 'e': 429, 'a': 430, 'm': 431, 'h': 432, 'd': 433, 'c': 434, 'g': 435, 'p': 436, 'f': 437, 'j': 438, 'v': 439, 'w': 440, 'x': 441, 'z': 442, 'q': 443, 's': 444, 't': 445, 'n': 446, 'r': 447, 'l': 448, 'o': 449, 'u': 450, 'i': 451, 'e': 452, 'a': 453, 'm': 454, 'h': 455, 'd': 456, 'c': 457, 'g': 458, 'p': 459, 'f': 460, 'j': 461, 'v': 462, 'w': 463, 'x': 464, 'z': 465, 'q': 466, 's': 467, 't': 468, 'n': 469, 'r': 470, 'l': 471, 'o': 472, 'u': 473, 'i': 474, 'e': 475, 'a': 476, 'm': 477, 'h': 478, 'd': 479, 'c': 480, 'g': 481, 'p': 482, 'f': 483, 'j': 484, 'v': 485, 'w': 486, 'x': 487, 'z': 488, 'q': 489, 's': 490, 't': 491, 'n': 492, 'r': 493, 'l': 494, 'o': 495, 'u': 496, 'i': 497, 'e': 498, 'a': 499, 'm': 500, 'h': 501, 'd': 502, 'c': 503, 'g': 504, 'p': 505, 'f': 506, 'j': 507, 'v': 508, 'w': 509, 'x': 510, 'z': 511, 'q': 512, 's': 513, 't': 514, 'n': 515, 'r': 516, 'l': 517, 'o': 518, 'u': 519, 'i': 520, 'e': 521, 'a': 522, 'm': 523, 'h': 524, 'd': 525, 'c': 526, 'g': 527, 'p': 528, 'f': 529, 'j': 530, 'v': 531, 'w': 532, 'x': 533, 'z': 534, 'q': 535, 's': 536, 't': 537, 'n': 538, 'r': 539, 'l': 540, 'o': 541, 'u': 542, 'i': 543, 'e': 544, 'a': 545, 'm': 546, 'h': 547, 'd': 548, 'c': 549, 'g': 550, 'p': 551, 'f': 552, 'j': 553, 'v': 554, 'w': 555, 'x': 556, 'z': 557, 'q': 558, 's': 559, 't': 560, 'n': 561, 'r': 562, 'l': 563, 'o': 564, 'u': 565, 'i': 566, 'e': 567, 'a': 568, 'm': 569, 'h': 570, 'd': 571, 'c': 572, 'g': 573, 'p': 574, 'f': 575, 'j': 576, 'v': 577, 'w': 578, 'x': 579, 'z': 580, 'q': 581, 's': 582, 't': 583, 'n': 584, 'r': 585, 'l': 586, 'o': 587, 'u': 588, 'i': 589, 'e': 590, 'a': 591, 'm': 592, 'h': 593, 'd': 594, 'c': 595, 'g': 596, 'p': 597, 'f': 598, 'j': 599, 'v': 600, 'w': 601, 'x': 602, 'z': 603, 'q': 604, 's': 605, 't': 606, 'n': 607, 'r': 608, 'l': 609, 'o': 610, 'u': 611, 'i': 612, 'e': 613, 'a': 614, 'm': 615, 'h': 616, 'd': 617, 'c': 618, 'g': 619, 'p': 620, 'f': 621, 'j': 622, 'v': 623, 'w': 624, 'x': 625, 'z': 626, 'q': 627, 's': 628, 't': 629, 'n': 630, 'r': 631, 'l': 632, 'o': 633, 'u': 634, 'i': 635, 'e': 636, 'a': 637, 'm': 638, 'h': 639, 'd': 640, 'c': 641, 'g': 642, 'p': 643, 'f': 644, 'j': 645, 'v': 646, 'w': 647, 'x': 648, 'z': 649, 'q': 650, 's': 651, 't': 652, 'n': 653, 'r': 654, 'l': 655, 'o': 656, 'u': 657, 'i': 658, 'e': 659, 'a': 660, 'm': 661, 'h': 662, 'd': 663, 'c': 664, 'g': 665, 'p': 666, 'f': 667, 'j': 668, 'v': 669, 'w': 670, 'x': 671, 'z': 672, 'q': 673, 's': 674, 't': 675, 'n': 676, 'r': 677, 'l': 678, 'o': 679, 'u': 680, 'i': 681, 'e': 682, 'a': 683, 'm': 684, 'h': 685, 'd': 686, 'c': 687, 'g': 688, 'p': 689, 'f': 690, 'j': 691, 'v': 692, 'w': 693, 'x': 694, 'z': 695, 'q': 696, 's': 697, 't': 698, 'n': 699, 'r': 700, 'l': 701, 'o': 702, 'u': 703, 'i': 704, 'e': 705, 'a': 706, 'm': 707, 'h': 708, 'd': 709, 'c': 710, 'g': 711, 'p': 712, 'f': 713, 'j': 714, 'v': 715, 'w': 716, 'x': 717, 'z': 718, 'q': 719, 's': 720, 't': 721, 'n': 722, 'r': 723, 'l': 724, 'o': 725, 'u': 726, 'i': 727, 'e': 728, 'a': 729, 'm': 730, 'h': 731, 'd': 732, 'c': 733, 'g': 734, 'p': 735, 'f': 736, 'j': 737, 'v': 738, 'w': 739, 'x': 740, 'z': 741, 'q': 742, 's': 743, 't': 744, 'n': 745, 'r': 746, 'l': 747, 'o': 748, 'u': 749, 'i': 750, 'e': 751, 'a': 752, 'm': 753, 'h': 754, 'd': 755, 'c': 756, 'g': 757, 'p': 758, 'f': 759, 'j': 760, 'v': 761, 'w': 762, 'x': 763, 'z': 764, 'q': 765, 's': 766, 't': 767, 'n': 768, 'r': 769, 'l': 770, 'o': 771, 'u': 772, 'i': 773, 'e': 774, 'a': 775, 'm': 776, 'h': 777, 'd': 778, 'c': 779, 'g': 780, 'p': 781, 'f': 782, 'j': 783, 'v': 784, 'w': 785, 'x': 786, 'z': 787, 'q': 788, 's': 789, 't': 790, 'n': 791, 'r': 792, 'l': 793, 'o': 794, 'u': 795, 'i': 796, 'e': 797, 'a': 798, 'm': 799, 'h': 800, 'd': 801, 'c': 802, 'g': 803, 'p': 804, 'f': 805, 'j': 806, 'v': 807, 'w': 808, 'x': 809, 'z': 8010, 'q': 8011, 's': 8012, 't': 8013, 'n': 8014, 'r': 8015, 'l': 8016, 'o': 8017, 'u': 8018, 'i': 8019, 'e': 8020, 'a': 8021, 'm': 8022, 'h': 8023, 'd': 8024, 'c': 8025, 'g': 8026, 'p': 8027, 'f': 8028, 'j': 8029, 'v': 8030, 'w': 8031, 'x': 8032, 'z': 8033, 'q': 8034, 's': 8035, 't': 8036, 'n': 8037, 'r': 8038, 'l': 8039, 'o': 8040, 'u': 8041, 'i': 8042, 'e': 8043, 'a': 8044, 'm': 8045, 'h': 8046, 'd': 8047, 'c': 8048, 'g': 8049, 'p': 8050, 'f': 8051, 'j': 8052, 'v': 8053, 'w': 8054, 'x': 8055, 'z': 8056, 'q': 8057, 's': 8058, 't': 8059, 'n': 8060, 'r': 8061, 'l': 8062, 'o': 8063, 'u': 8064, 'i': 8065, 'e': 8066, 'a': 8067, 'm': 8068, 'h': 8069, 'd': 8070, 'c': 8071, 'g': 8072, 'p': 8073, 'f': 8074, 'j': 8075, 'v': 8076, 'w': 8077, 'x': 8078, 'z': 8079, 'q': 8080, 's': 8081, 't': 8082, 'n': 8083, 'r': 8084, 'l': 8085, 'o': 8086, 'u': 8087, 'i': 8088, 'e': 8089, 'a': 8090, 'm': 8091, 'h': 8092, 'd': 8093, 'c': 8094, 'g': 8095, 'p': 8096, 'f': 8097, 'j': 8098, 'v': 8099, 'w': 80100, 'x': 80101, 'z': 80102, 'q': 80103, 's': 80104, 't': 80105, 'n': 80106, 'r': 80107, 'l': 80108, 'o': 80109, 'u': 80110, 'i': 80111, 'e': 80112, 'a': 80113, 'm': 80114, 'h': 80115, 'd': 80116, 'c': 80117, 'g': 80118, 'p': 80119, 'f': 80120, 'j': 80121, 'v': 80122, 'w': 80123, 'x': 80124, 'z': 80125, 'q': 80126, 's': 80127, 't': 80128, 'n': 80129, 'r': 80130, 'l': 80131, 'o': 80132, 'u': 80133, 'i': 80134, 'e': 80135, 'a': 80136, 'm': 80137, 'h': 80138, 'd': 80139, 'c': 80140, 'g': 80141, 'p': 80142, 'f': 80143, 'j': 80144, 'v': 80145, 'w': 80146, 'x': 80147, 'z': 80148, 'q': 80149, 's': 80150, 't': 80151, 'n': 80152, 'r': 80153, 'l': 80154, 'o': 80155, 'u': 80156, 'i': 80157, 'e': 80158, 'a': 80159, 'm': 80160, 'h': 80161, 'd': 80162, 'c': 80163, 'g': 80164, 'p': 80165, 'f': 80166, 'j': 80167, 'v': 80168, 'w': 80169, 'x': 80170, 'z': 80171, 'q': 80172, 's': 80173, 't': 80174, 'n': 80175, 'r': 80176, 'l': 80177, 'o': 80178, 'u': 80179, 'i': 80180, 'e': 80181, 'a': 80182, 'm': 80183, 'h': 80184, 'd': 80185, 'c': 80186, 'g': 80187, 'p': 80188, 'f': 80189, 'j': 80190, 'v': 80191, 'w': 80192, 'x': 80193, 'z': 80194, 'q': 80195, 's': 80196, 't': 80197, 'n': 80198, 'r': 80199, 'l': 80200, 'o': 80201, 'u': 80202, 'i': 80203, 'e': 80204, 'a': 80205, 'm': 80206, 'h': 80207, 'd': 80208, 'c': 80209, 'g': 80210, 'p': 80211, 'f': 80212, 'j': 80213, 'v': 80214, 'w': 80215, 'x': 80216, 'z': 80217, 'q': 80218, 's': 80219, 't': 80220, 'n': 80221, 'r': 80222, 'l': 80223, 'o': 80224, 'u': 80225, 'i': 80226, 'e': 80227, 'a': 80228, 'm': 80229, 'h': 80230, 'd': 80231, 'c': 80232, 'g': 80233, 'p': 80234, 'f': 80235, 'j': 80236, 'v': 80237, 'w': 80238, 'x': 80239, 'z': 80240, 'q': 80241, 's': 80242, 't': 80243, 'n': 80244, 'r': 80245, 'l': 80246, 'o': 80247, 'u': 80248, 'i': 80249, 'e': 80250, 'a': 80251, 'm': 80252, 'h': 80253, 'd': 80254, 'c': 80255, 'g': 80256, 'p': 80257, 'f': 80258, 'j': 80259, 'v': 80260, 'w': 80261, 'x': 80262, 'z': 80263, 'q': 80264, 's': 80265, 't': 80266, 'n': 80267, 'r': 80268, 'l': 80269, 'o': 80270, 'u': 80271, 'i': 80272, 'e': 80273, 'a': 80274, 'm': 80275, 'h': 80276, 'd': 80277, 'c': 80278, 'g': 80279, 'p': 80280, 'f': 80281, 'j': 80282, 'v': 80283, 'w': 80284, 'x': 80285, 'z': 80286, 'q': 80287, 's': 80288, 't': 80289, 'n': 80290, 'r': 80291, 'l': 80292, 'o': 80293, 'u': 80294, 'i': 80295, 'e': 80296, 'a': 80297, 'm': 80298, 'h': 80299, 'd': 80300, 'c': 80301, 'g': 80302, 'p': 80303, 'f': 80304, 'j': 80305, 'v': 80306, 'w': 80307, 'x': 80308, 'z': 80309, 'q': 80310, 's': 80311, 't': 80312, 'n': 80313, 'r': 80314, 'l': 80315, 'o': 80316, 'u': 80317, 'i': 80318, 'e': 80319, 'a': 80320, 'm': 80321, 'h': 80322, 'd': 80323, 'c': 80324, 'g': 80325, 'p': 80326, 'f': 80327, 'j': 80328, 'v': 80329, 'w': 80330, 'x': 80331, 'z': 80332, 'q': 80333, 's': 80334, 't': 80335, 'n': 80336, 'r': 80337, 'l': 80338, 'o': 80339, 'u': 80340, 'i': 80341, 'e': 80342, 'a': 80343, 'm': 80344, 'h': 80345, 'd': 80346, 'c': 80347, 'g': 80348, 'p': 80349, 'f': 80350, 'j': 80351, 'v': 80352, 'w': 80353, 'x': 80354, 'z': 80355, 'q': 80356, 's': 80357, 't': 80358, 'n': 80359, 'r': 80360, 'l': 80361, 'o': 80362, 'u': 80363, 'i': 80364, 'e': 80365, 'a': 80366, 'm': 80367, 'h': 80368, 'd': 80369, 'c': 80370, 'g': 80371, 'p': 80372, 'f': 80373, 'j': 80374, 'v': 80375, 'w': 80376, 'x': 80377, 'z': 80378, 'q': 80379, 's': 80380, 't': 80381, 'n': 80382, 'r': 80383, 'l': 80384, 'o': 80385, 'u': 80386, 'i': 80387, 'e': 80388, 'a': 80389, 'm': 80390, 'h': 80391, 'd': 80392, 'c': 80393, 'g': 80394, 'p': 80395, 'f': 80396, 'j': 80397, 'v': 80398, 'w': 80399, 'x': 80400, 'z': 80401, 'q': 80402, 's': 80403, 't': 80404, 'n': 80405, 'r': 80406, 'l': 80407, 'o': 80408, 'u': 80409, 'i': 80410, 'e': 80411, 'a': 80412, 'm': 80413, 'h': 80414, 'd': 80415, 'c': 80416, 'g': 80417, 'p': 80418, 'f': 80419, 'j': 80420, 'v': 80421, 'w': 80422, 'x': 80423, 'z': 80424, 'q': 80425, 's': 80426, 't': 80427, 'n': 80428, 'r': 80429, 'l': 80430, 'o': 80431, 'u': 80432, 'i': 80433, 'e': 80434, 'a': 80435, 'm': 80436, 'h': 80437, 'd': 80438, 'c': 80439, 'g': 80440, 'p': 80441, 'f': 80442, 'j': 80443, 'v': 80444, 'w': 80445, 'x': 80446, 'z': 80447, 'q': 80448, 's': 80449, 't': 80450, 'n': 80451, 'r': 80452, 'l': 80453, 'o': 80454, 'u': 80455, 'i': 80456, 'e': 80457, 'a': 80458, 'm': 80459, 'h': 80460, 'd': 80461, 'c': 80462, 'g': 80463, 'p': 80464, 'f': 80465, 'j': 80466, 'v': 80467, 'w': 80468, 'x': 80469, 'z': 80470, 'q': 80471, 's': 80472, 't': 80473, 'n': 80474, 'r': 80475, 'l': 80476, 'o': 80477, 'u': 80478, 'i': 80479, 'e': 80480, 'a': 80481, 'm': 80482, 'h': 80483, 'd': 80484, 'c': 80485, 'g': 80486, 'p': 80487, 'f': 80488, 'j': 80489, 'v': 80490, 'w': 80491, 'x': 80492, 'z': 80493, 'q': 80494, 's': 80495, 't': 80496, 'n': 80497, 'r': 80498, 'l': 80499, 'o': 80500, 'u': 80501, 'i': 80502, 'e': 805
```

```
#Creating a character embedding matrix
embd_ch_matrix=np.zeros((41,41))
for x,y in ke_tokenizer2.word_index.items():
    embd_ch_matrix[y][y]=1
```

```
embd_ch_matrix
```

```
array([[0., 0., 0., ..., 0., 0., 0.],
       [0., 1., 0., ..., 0., 0., 0.],
       [0., 0., 1., ..., 0., 0., 0.],
       ...,
       [0., 0., 0., ..., 1., 0., 0.],
       [0., 0., 0., ..., 0., 1., 0.],
       [0., 0., 0., ..., 0., 0., 1.]])
```

```
embd_ch_matrix.shape
```

```
(41, 41)
```

```
#Model 2 Layers def
#Defining Embedding Layer
Embd_ly_char=Embedding(len(ke_tokenizer2.word_index)+1,41,embeddings_initializer=tf.keras.initi

input_layer2=Input(shape=(max_ch_len,))
embd_layer2=Embd_ly_char(input_layer2)
N1=Conv1D(64,3,activation='relu',kernel_initializer=tf.keras.initializers.he_normal(),kernel_
M1=Conv1D(64,3,activation='relu',kernel_initializer=tf.keras.initializers.he_normal(),kernel_

Max_pool1=MaxPool1D(5)(M1)
K1=Conv1D(64,3,activation='relu',kernel_initializer=tf.keras.initializers.he_normal(),kernel_
T1=Conv1D(64,3,activation='relu',kernel_initializer=tf.keras.initializers.he_normal(),kernel_
Max_pool2=MaxPool1D(4)(T1)
flat_layer1=Flatten()(Max_pool2)
dropout_layer1=Dropout(0.3)(flat_layer1)
dense_layer1=Dense(64,activation='relu',kernel_initializer=tf.keras.initializers.he_normal(),
Output_layer1=Dense(20,activation='softmax',kernel_initializer=tf.keras.initializers.glorot_n
M2=Model(inputs=input_layer2,outputs=Output_layer1)
```

```
M2.summary()
```

```
Model: "model_3"
```

Layer (type)	Output Shape	Param #
<hr/>		
input_4 (InputLayer)	[(None, 8062)]	0
embedding_2 (Embedding)	(None, 8062, 41)	1681
conv1d_21 (Conv1D)	(None, 8060, 64)	7936

conv1d_22 (Conv1D)	(None, 8058, 64)	12352
max_pooling1d_6 (MaxPooling1D)	(None, 1611, 64)	0
conv1d_23 (Conv1D)	(None, 1609, 64)	12352
conv1d_24 (Conv1D)	(None, 1607, 64)	12352
max_pooling1d_7 (MaxPooling1D)	(None, 401, 64)	0
flatten_3 (Flatten)	(None, 25664)	0
dropout_3 (Dropout)	(None, 25664)	0
dense_6 (Dense)	(None, 64)	1642560
dense_7 (Dense)	(None, 20)	1300
<hr/>		
Total params: 1,690,533		
Trainable params: 1,688,852		
Non-trainable params: 1,681		

```

## f1_score_callback
f1_metric=f1((X_test,Y_test))

## Callback for saving best model
checkpoint = ModelCheckpoint(filepath='best_model_1.h5',verbose=1,monitor='val_accuracy',
                             mode='max',save_best_only=True)

## Callback for earlystopping
early_stop = EarlyStopping(monitor="val_accuracy",mode='max',patience=2)

## Tensorboard
log_dir = "logs"
tensorboard_cb = TensorBoard(log_dir=log_dir,histogram_freq=1,write_graph=True)

#reduceon Plateau
reduce_cb=ReduceLROnPlateau(monitor='val_accuracy',factor=0.9,patience=1,verbose=1,mode='max')

## all callbacks
callbacks =[f1_metric,early_stop,tensorboard_cb,checkpoint,reduce_cb]

## compile model
M2.compile(loss='categorical_crossentropy', optimizer=Adam(learning_rate=0.001), metrics=['ac'])

## Trainning

M2.fit(X_train,Y_train,epochs=20,verbose=2,validation_data=(X_test,Y_test),batch_size =128,ca

Epoch 1/20

```

```
111/111 - 240s - loss: 35.5907 - accuracy: 0.0554 - val_loss: 13.0939 - val_accuracy: 0
Validation F1 score 0.052899936265137025

Epoch 00001: val_accuracy improved from -inf to 0.05290, saving model to best_model_1.h5
Epoch 2/20
111/111 - 241s - loss: 8.0364 - accuracy: 0.0513 - val_loss: 5.4160 - val_accuracy: 0.053112385808370514
Validation F1 score 0.052899936265137025

Epoch 00002: val_accuracy did not improve from 0.05290

Epoch 00002: ReduceLROnPlateau reducing learning rate to 0.0009000000427477062.
Epoch 3/20
111/111 - 239s - loss: 4.8527 - accuracy: 0.0514 - val_loss: 4.7086 - val_accuracy: 0.053112385808370514
Validation F1 score 0.053112385808370514

Epoch 00003: val_accuracy improved from 0.05290 to 0.05311, saving model to best_model_1.h5
Epoch 4/20
111/111 - 249s - loss: 4.7344 - accuracy: 0.0530 - val_loss: 4.7371 - val_accuracy: 0.053112385808370514
Validation F1 score 0.053112385808370514

Epoch 00004: val_accuracy did not improve from 0.05311

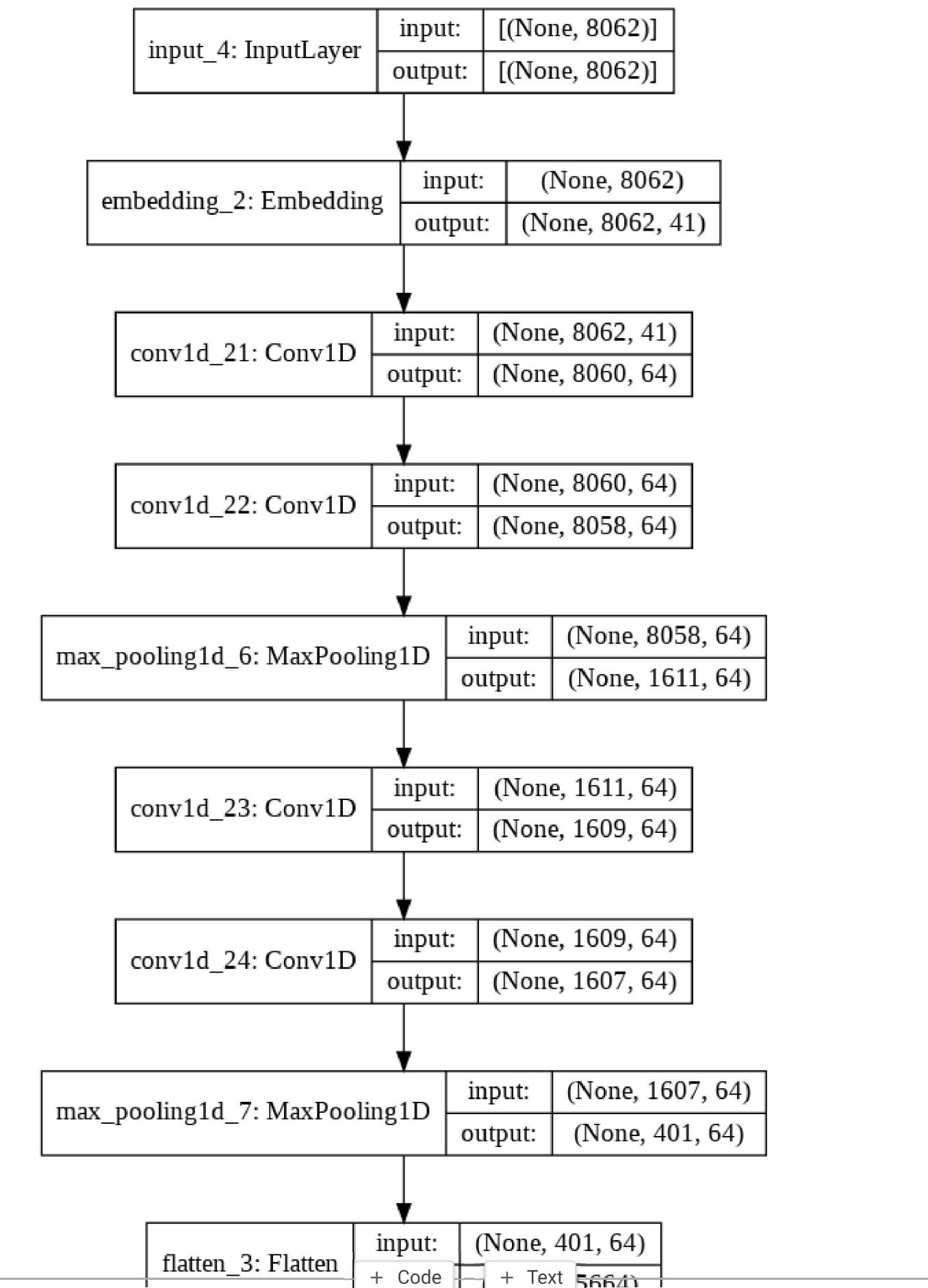
Epoch 00004: ReduceLROnPlateau reducing learning rate to 0.0008100000384729356.
Epoch 5/20
111/111 - 247s - loss: 4.5389 - accuracy: 0.0530 - val_loss: 4.5227 - val_accuracy: 0.053112385808370514
Validation F1 score 0.053112385808370514

Epoch 00005: val_accuracy did not improve from 0.05311

Epoch 00005: ReduceLROnPlateau reducing learning rate to 0.0007290000503417104.
<keras.callbacks.History at 0x7fb6571bb1d0>
```

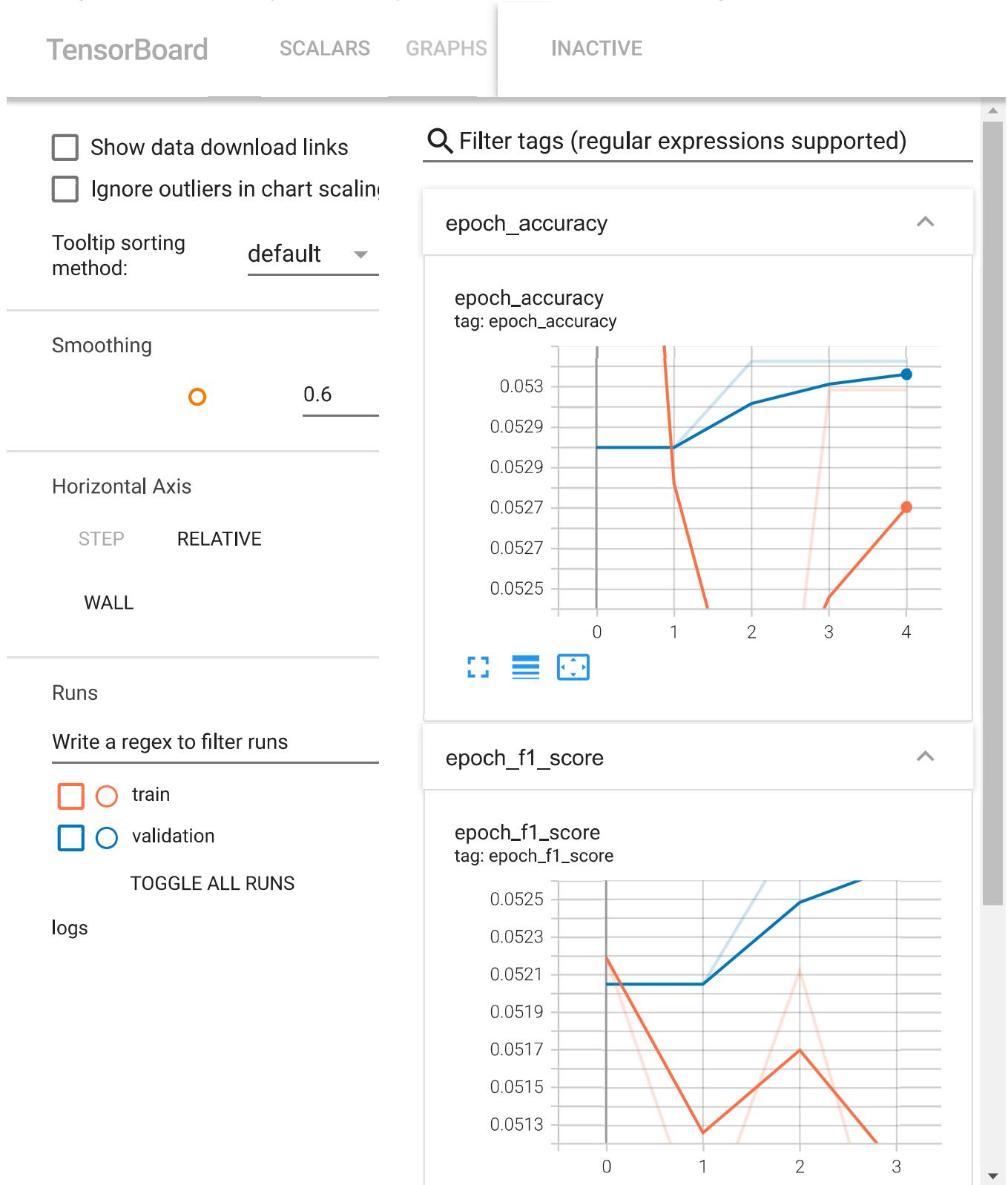
```
tf.keras.utils.plot_model(M2, to_file='M2.png', show_shapes=True, show_layer_names=True)
```

→



```
%tensorboard --logdir logs
```

Reusing TensorBoard on port 6006 (pid 5878), started 0:38:49 ago. (Use '!kill 5878' to kill)



✓ 0s completed at 11:48 PM

